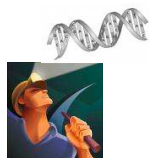


# Mineração de Dados em Biologia Molecular

## Agrupamento de Dados

André C. P. L. F. de Carvalho  
Monitor: Valéria Carvalho



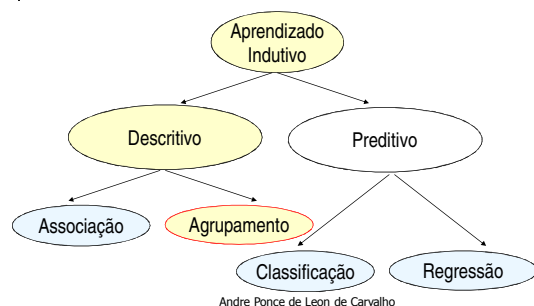
## Tópicos

- Agrupamento de dados
- Análise de cluster
- Dificuldades em agrupamento
- Algoritmos de agrupamento
- Validação
- Aplicações

Andre Ponce de Leon de Carvalho

2

## Agrupamento



Andre Ponce de Leon de Carvalho

3

## Agrupamento

- Organização de um conjunto de objetos em grupos (clusters)
  - De acordo com alguma forma de similaridade ou relação entre eles

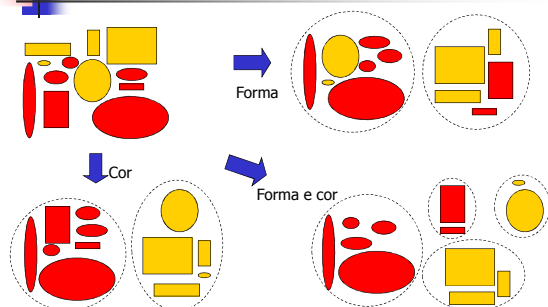


Como organizar?

Andre Ponce de Leon de Carvalho

4

## Agrupamento



Andre Ponce de Leon de Carvalho

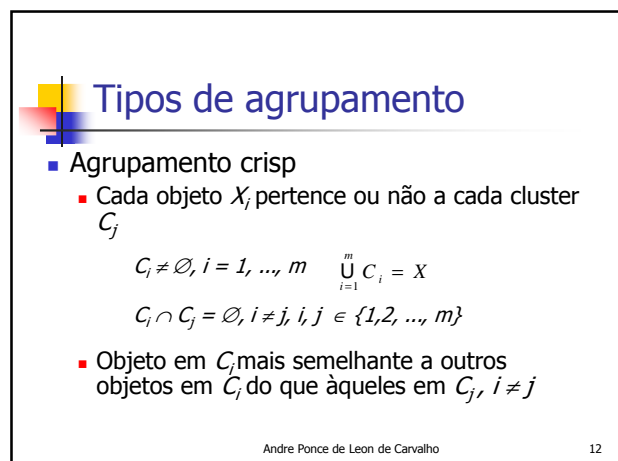
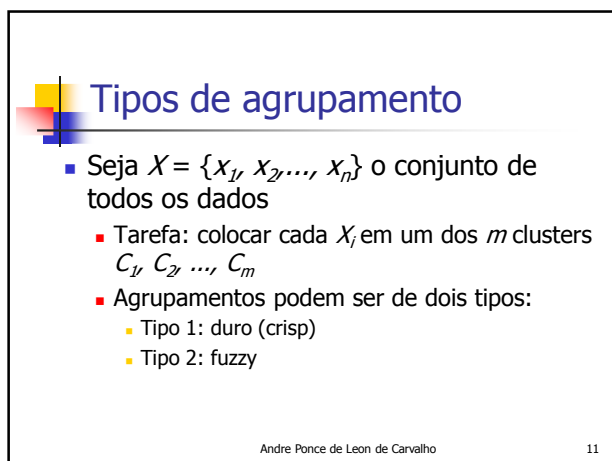
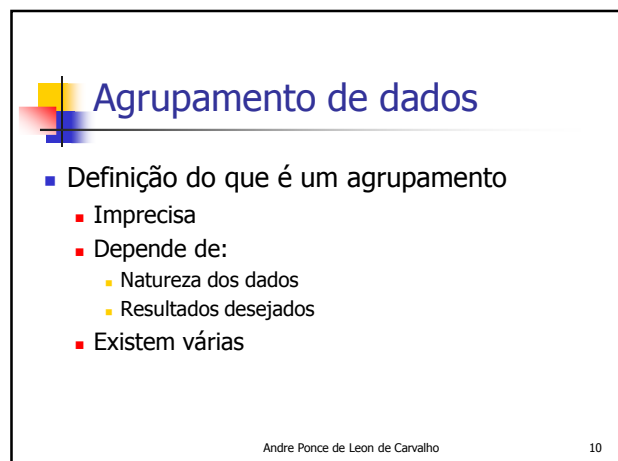
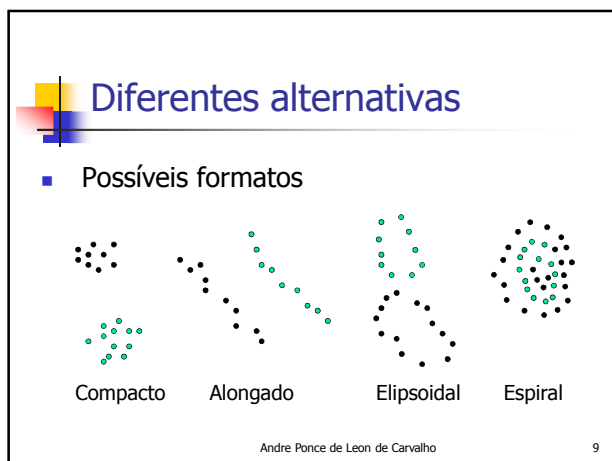
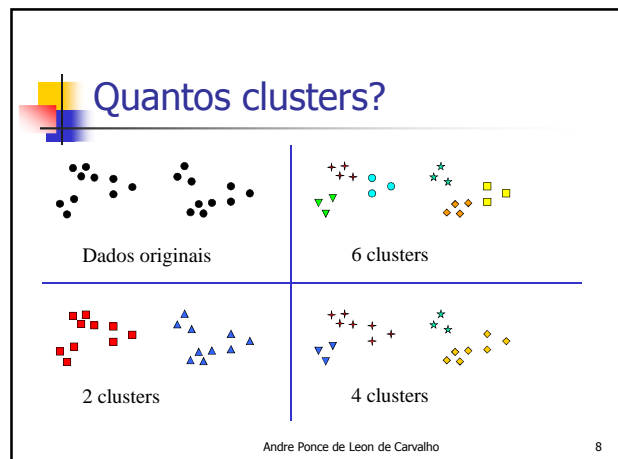
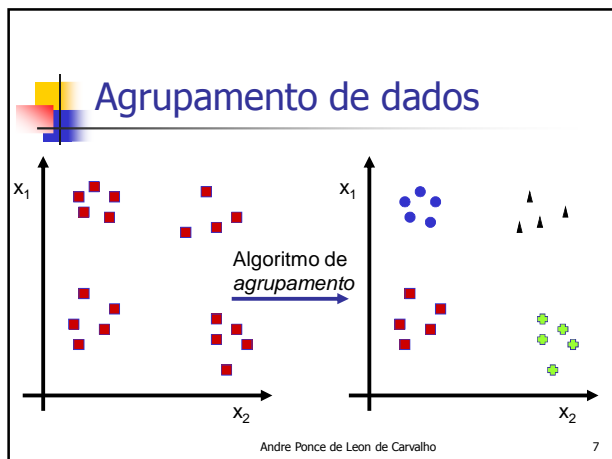
5

## Objetivo

- Encontrar conjunto de clusters que maximizam SM ou minimizam DM
  - SM: medida de similaridade
  - DM: medida de dissimilaridade
  - Quanto maior a homogeneidade dentro dos grupos e a diferença entre os grupos, melhor
- Alternativas
  - Busca exaustiva
  - Técnicas mais sofisticadas

Andre Ponce de Leon de Carvalho

6



## Tipos de agrupamento

### ■ Agrupamento fuzzy

- Usa uma função de pertinência para definir o quanto um elemento pertence a um grupo

$$\mu_j : X \rightarrow [0, 1]$$

$$\sum_{j=1}^m \mu_j(x_i) = 1, i \in \{1, \dots, n\} \quad \begin{array}{l} m = \text{número de grupos} \\ n = \text{número de objetos} \end{array}$$

$$0 < \sum_{i=1}^n \mu_j(x_i) < n, j \in \{1, \dots, m\}$$

Andre Ponce de Leon de Carvalho

13

## Algoritmos de agrupamento

### ■ Busca exaustiva

- Tentar todos os possíveis agrupamentos de tamanho  $m$  (para vários valores de  $m$ )
- Números de Stirling do segundo tipo
  - Número de formas de particionar  $n$  dados em  $m$  subconjuntos não vazios

$$>> \binom{n}{m} \geq \left(\frac{n}{m}\right)^m \quad \begin{array}{l} m = \text{número de grupos} \\ n = \text{número de objetos} \end{array}$$

- Impraticável

Andre Ponce de Leon de Carvalho

14

## Algoritmos de agrupamento

### ■ Particionais

- Protótipos (erro quadrático médio)
- Densidade

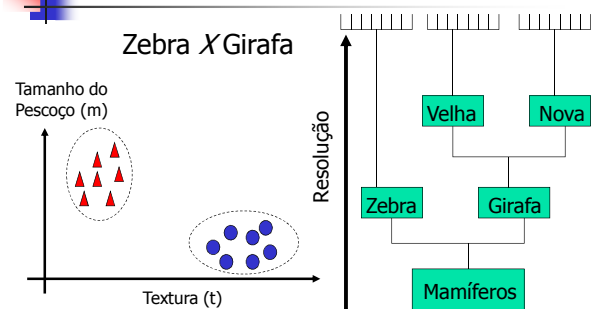
### ■ Hierárquicos

- Baseados em grids
- Baseados em grafos
- Outros algoritmos
  - Ex.: Redes neurais SOM

Andre Ponce de Leon de Carvalho

15

## Particional X Hierárquico



Andre Ponce de Leon de Carvalho

16

## Algoritmos particionais

### ■ Principais características

- Produzem um único agrupamento
- A maioria utiliza abordagem "gulosa" (*greedy*)
  - Sempre procura escolher melhor alternativa atual, sem considerar consequências futuras
  - Uma vez tomada uma decisão, ela não é mais alterada
  - Geralmente resultado depende da ordem de apresentação dos exemplos

Andre Ponce de Leon de Carvalho

17

## Algoritmos particionais

- K-médias
- K-médias ótimo
- K-médias seqüencial
- SOM
- FCM
- DENCLUE
- CLICK
- CAST
- SNN

Andre Ponce de Leon de Carvalho

18

## Algoritmo k-médias

- Supor  $n$  objetos  $x_1, x_2, \dots, x_n$  a serem agrupados em  $k$  clusters,  $k < n$ 
  - Seja  $\mu_i$  a média dos objetos do cluster  $C_i$
  - Medida de distância pode ser utilizada para definir a que cluster um objeto pertence
    - $x_p \in \text{cluster } C_i$  se  $d(x_p, \mu_i)$  é menor que todas as  $k-1$  distâncias entre  $x_p$  e  $\mu_j$ ,  $j = 1, 2, \dots, k$  e  $i \neq j$

Andre Ponce de Leon de Carvalho

19

## Algoritmo k-médias

```
1 Sugerir médias  $\mu_1, \mu_2, \dots, \mu_k$  iniciais
2 Repetir
    Usar as médias sugeridas para agrupar
    os objetos em  $K$  clusters
    Para  $i$  variando de 1 a  $K$ 
        Substituir  $\mu_i$  pela média de todos os
        objetos do cluster  $C_i$ 
    Até nenhuma das médias mudar
```

Andre Ponce de Leon de Carvalho

20

## Algoritmo k-médias

- Médias iniciais
  - Exemplos (vetores) aleatórios
  - Elementos aleatoriamente escolhidos do conjunto de treinamento

Andre Ponce de Leon de Carvalho

21

## Limitações do k-médias

- Escolha do valor de  $K$
- K-médias tem problemas quando os grupos têm:
  - Diferentes densidades
  - Formatos não hiper-esféricos
- Tem problemas também quando os dados contêm *outliers*

Andre Ponce de Leon de Carvalho

22

## Exercício

- Agrupar, utilizando k-médias, os dados abaixo em 2 grupos:
  - $X_1 = 1, 0, 1, 1$
  - $X_2 = 1, 0, 0, 0$
  - $X_3 = 0, 1, 1, 0$
  - $X_4 = 1, 1, 1, 1$
  - $X_5 = 0, 0, 0, 1$

Andre Ponce de Leon de Carvalho

23

## Validação de agrupamentos

- Existem várias medidas para avaliar qualidade de classificadores
  - Acurácia, precisão, revocação, F1
- Como avaliar os clusters gerados por um algoritmo de agrupamento?

Andre Ponce de Leon de Carvalho

24

## Medidas de validação

- Existem várias medidas de validação
  - Julgam aspectos diferentes
- Podem ser divididas em três grupos
  - Índices ou critérios externos
    - Medem o quanto os rótulos dos grupos casam com a classe verdadeira
  - Índices ou critérios internos
    - Medem a qualidade da partição obtida sem considerar informações externas
  - Índices ou critérios relativos
    - Usados para comparar duas partições ou grupos

Andre Ponce de Leon de Carvalho

25

## Medidas internas

- Coesão de clusters
  - Mede o quão relacionados estão os objetos dentro de um cluster
- Separação de clusters
  - Mede o quão distinto ou separado cada cluster é dos demais clusters
- Silhueta

Andre Ponce de Leon de Carvalho

26

## Silhueta

- Combina coesão com separação
- Calculada para cada objeto que faz parte de um agrupamento
  - Baseada na proximidade entre os objetos de um cluster e na distância dos objetos de um cluster ao cluster mais próximo

Andre Ponce de Leon de Carvalho

27

## Silhueta

- Para cada objeto  $x_i$ 
    - $a(x_i)$ : distância média de  $x_i$  aos outros objetos de seu cluster
    - $b(x_i)$ : min (distância média de  $x_i$  a todos os objetos de cada outro cluster)
- $$s(x_i) = \begin{cases} 1 - a(x_i) / b(x_i), & \text{se } a(x_i) < b(x_i) \\ 0, & \text{se } a(x_i) = b(x_i) \\ b(x_i) / a(x_i) - 1, & \text{se } a(x_i) > b(x_i) \end{cases}$$
- Largura média da silhueta
    - Média sobre todos os objetos do conjunto de dados
    - Valor entre -1 e 1 (quanto mais próximo de 1, melhor)

Andre Ponce de Leon de Carvalho

28

## Algoritmos hierárquicos

- Utilizam diagrama de árvore (dendograma)
  - Produzem uma sequência (hierarquia) de agrupamentos
- Historicamente utilizados em áreas que utilizam estrutura hierárquica de dados
  - Ex.: Biologia e arqueologia

Andre Ponce de Leon de Carvalho

29

## Algoritmos hierárquicos

- Conceito de representação hierárquica de dados originou-se na Biologia
  - Algoritmos de agrupamento hierárquicos  $\equiv$  estrutura hierárquica da taxonomia de Linnaean
  - Biólogos geralmente preferem agrupamentos hierárquicos

Andre Ponce de Leon de Carvalho

30

## Algoritmos hierárquicos

- Aplicações na biologia geralmente não se preocupam com o número ótimo de clusters
  - Biólogos geralmente estão interessados na hierarquia completa

Andre Ponce de Leon de Carvalho

31

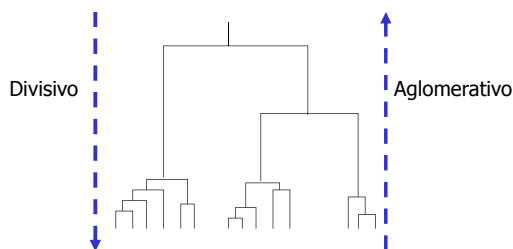
## Algoritmos hierárquicos

- Podem ser de dois tipos:
  - Aglomerativos: combinam, repetidamente, dois clusters em um
    - A cada passo, combina os dois clusters mais próximos
  - Divisivos: Dividem, repetidamente, um cluster em dois
    - A cada passo, divide o cluster menos homogêneo em dois novos clusters

Andre Ponce de Leon de Carvalho

32

## Exemplo



Andre Ponce de Leon de Carvalho

33

## Algoritmos aglomerativos

- Começam com  $C_0 = \{\{x_1\}, \dots, \{x_n\}\}$
- A cada passo  $t$ , combinam dois clusters em um, produzindo:
  - $|C_{t+1}| = |C_t| - 1$  e  $C_t \subset C_{t+1}$
- No passo final (passo  $n-1$ ) tem-se a hierarquia:
  - $C_0 = \{\{x_1\}, \dots, \{x_n\}\} \subset C_1 \dots \subset C_{n-1} = \{x_1, \dots, x_n\}$

Andre Ponce de Leon de Carvalho

34

## Algoritmos divisivos

- Começam com  $C_0 = \{x_1, \dots, x_n\}$
- A cada passo  $t$ , dividem um cluster em dois, produzindo:
  - $|C_{t+1}| = |C_t| + 1$  e  $C_{t+1} \subset C_t$
- No passo final (passo  $n-1$ ) tem-se a hierarquia:
  - $C_{n-1} = \{\{x_1\}, \dots, \{x_n\}\} \subset \dots \subset C_0 = \{x_1, \dots, x_n\}$

Andre Ponce de Leon de Carvalho

35

## Algoritmo aglomerativo

```

1 Inicializar  $C_0 = \{\{x_1\}, \dots, \{x_n\}\}$ 
2 Para  $t = 1$  até  $n - 1$  faça
    Encontrar o par de clusters mais próximos ( $C_i, C_j$ )
     $C_t = (C_{t-1} - \{C_i, C_j\}) \cup \{C_i \cup C_j\}$  /* atualizar centros
    
```

Andre Ponce de Leon de Carvalho

36

## Algoritmos hierárquicos

- Existe uma grande variedade de algoritmos hierárquicos

- Geralmente diferem na forma de calcular distância inter-clusters

$$d_{AB} = \min_{\substack{i \in A \\ j \in B}} (d_{ij}) \quad \text{Por ligação simples (single-link)}$$

$$d_{AB} = \max_{\substack{i \in A \\ j \in B}} (d_{ij}) \quad \text{Por ligação completa (complete-link)}$$

$$d_{AB} = \frac{1}{n_A n_B} \sum_{i \in A} \sum_{j \in B} d_{ij} \quad \text{Pela média do grupo (average-link)}$$

Andre Ponce de Leon de Carvalho

37

## Algoritmos hierárquicos

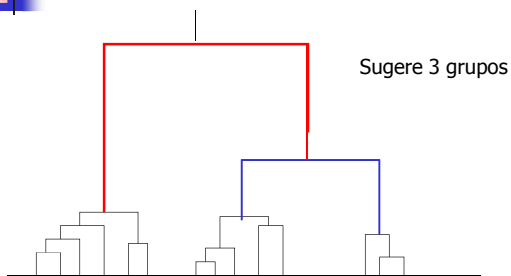
- Como escolher uma partição?

- Partição com  $n$  clusters
  - Selecionando partição com  $n$  clusters na sequência de agrupamentos da hierarquia
- Partição que melhor se encaixa nos dados
  - Procurar no dendograma grandes mudanças em níveis adjacentes
    - Nesse caso, uma mudança de  $j$  para  $j-1$  grupos pode indicar que  $j$  é o melhor número de grupos
    - Existem outros procedimentos, alguns mais objetivos

Andre Ponce de Leon de Carvalho

38

## Exemplo



Andre Ponce de Leon de Carvalho

39

## Exercício

- Seja o seguinte cadastro de pacientes:

Nome	Febre	Enjôo	Manchas	Dores	Diagnóstico
João	sim	sim	pequenas	sim	doente
Pedro	não	não	grandes	não	saudável
Maria	sim	sim	pequenas	não	saudável
José	sim	não	grandes	sim	doente
Ana	sim	não	pequenas	sim	saudável
Leila	não	não	grandes	sim	doente

Andre Ponce de Leon de Carvalho

40

## Exercício

- Agrupar os dados em dois grupos usando o algoritmo K-médias e medida de silhueta
  - Usar  $k = 2$  e  $k = 3$
  - Informação sobre a classe não deve ser usada
- Em que grupos seriam colocados os novos casos?
  - (Luis, não, não, pequenas, sim)
  - (Laura, sim, sim, grandes, sim)

Andre Ponce de Leon de Carvalho

41

## Considerações finais

- Abordagens tradicionais de agrupamento são muito utilizadas em AM
  - Várias definições de agrupamento
  - Diversos algoritmos
- Dificuldade de validar agrupamentos encontrados
- Semi-supervisionado

Andre Ponce de Leon de Carvalho

42



## Perguntas

---



Andre Ponce de Leon de Carvalho

43